



LAWRENCE  
LIVERMORE  
NATIONAL  
LABORATORY

# STK Bladestore Tests

Todd Heer

December 12, 2003

## **Disclaimer**

---

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This work was performed under the auspices of the U.S. Department of Energy by University of California, Lawrence Livermore National Laboratory under Contract W-7405-Eng-48.

## Overview

TheSTKBladestoreisadisksubsystemconsistingofATA disks,fiberchannel connectivity,andaRAIDcontroller(LSImanufactured).Thereareessentiallyfourhost connections\*andfourbackendfiberconnections.Thehostsideportsare2Gb/sec and withtheiradvertised400MB/secbandwidth,thedisksideportsare1GB/sec.

Ourgoalistotestthisflavorofdisktoseewhattherealworldperformancemightbe.

\*Actually8portsacrossfourcards,buttheyrecommendtouseonlyoneportone achcardasthoseare basicallymini -hubs

## Hardware

TheborrowedSTKBladestoreconsistedof2“B150”diskdrawerseachwith10blades offive250GBdrivesforatotalrawcapacityofjustunder25TBintwodrawers,one “BC84”ControlModule,andthe“F40 ”40Uheightcabinet.

In16Uofheightina40Ucabinet,therewas567lbsofhardware.Thatissomethingto makenoteof.Afullrackwouldweighinat1147pounds.

ThehostsystemisanIBM7026- 6M1with8GBofRAM,8PowerPC\_RS64IV processors,and fourseparateRIOdrawers.TheHBAsareIBM6228fiberchannelcards whichare2gigabitpersecond(gb/s)speedcapable.Thediskswereconnectednatively tothe6M1withouthouseofafiberswitch(suchasaBrocade).Thissystemhas sustainedI/Othroughputsofover600MB/sec,soitwasmorethanadequateforhistest. Itisalsorepresentativeofwhatweruninourproductionstorageenvironment.

TheoperatingsystemisAIX5.1withmaintenancelevel3applied.

## Limitations

Connectednatively aswewere,thereisanapparenttwo -pathlimitationfromthehost. ThismeansthatIwouldonlybeabletotesttwoHBAsworthofthroughput.AlthoughI haven'tverifiedit,mysuspicionisthatthiswouldbealleviatedinaSANenvironment - perhapsaBrocadefiberswitchin -betweenthedisksubsystemandthehost.IBMsupport verifiedthis.Givenmoretime,Iwouldtestthatassumption.

Itshouldalsobenotedthatittakesquiteabitoftimetoreconfigure/formatnewRAID groupsintheBladestore,as thediskisquitedenseforitssize.Wefirstconfigured2TB LUNsandthattookover18hourstoformat.

## Tests

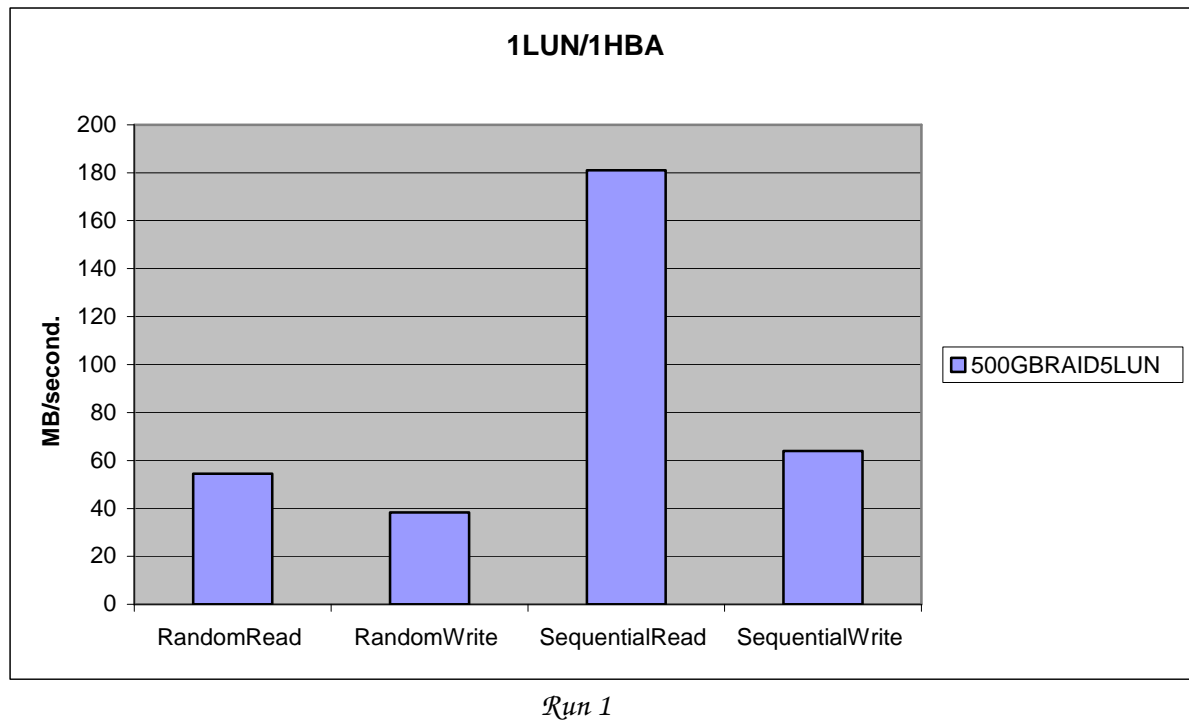
Allrunshave2typesoftestsrunagainstthem:randomaccessandsequential.The randomaccessteststakeaLUN,splitsitinto16llogicalvolumes,andthenhammers away(readsorwritesdependingonwhichpartofthetestisbeingrun)with4MBbuffers atanypointinthelogicalvolumefor10minutes(5minutesforwritesthen5minutesfor reads).Itisknowas“donnie.”

Theseq uentialtestissimplytheddcommandwith16MBbuffers.Writesaresimplya stringofzero’s(/dev/zero).TheystartfromthebeginningoftheLUNand/orlogical volume(numbershaveproventobethesameeitherway)andareallowedtorunfora periodof timesufficienttogarnergoodnumbers.

## TheHardFacts

Alltestswererunagainst500GBLUNs.ThisistheLUNsizethatseemedtoperformthe best,wasabletoberecognizedbyAIX(2TBLUNsweren’tabletobeconfigured),and offerenoughsizetobep otentiallyusefulwithoutbeingtoolargeforourintended purpose.Thesegmentsizeonthecontrollerwasobservedtobe128KB.

ThefirstrunwasasingleLUNoverasinglepath.Onlyonediskisaccessed,making useofoneHBA.Asthereisnocontenti onbetweendisksorHBAs,this test’sresultsare quiteusefulasabaselineandaiddindeterminingwherebottlenecks layinfuturetests.



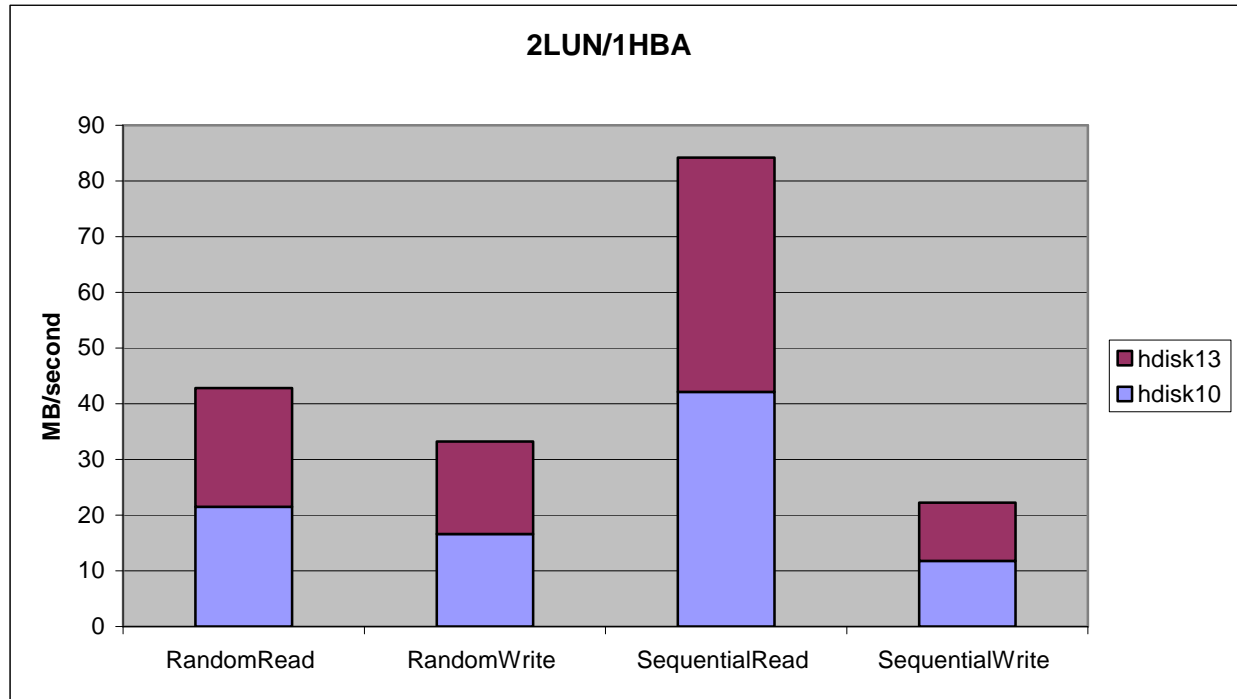
Inthistestweseethattherandomreadsarearound54MB/secandrandomwrites are 38MB/sec.Sequentialreadsareimpressiveat181MB/sec,whichisnearthetheoretical limitof2GB/seclinespeeds.Sequentialwritesarelowat64MB/sec.Itwasobservedthat

## STKBladestoreTests

August-September2003

thefirstfewseconds( $\leq 4$ )ofawritetestwereapproximately50%fastert hananyother timeintheremainderofthatspecificrun.Isuspectthisisareultofcachingandisreally notatruemeasureofsustainedwriteperformance.

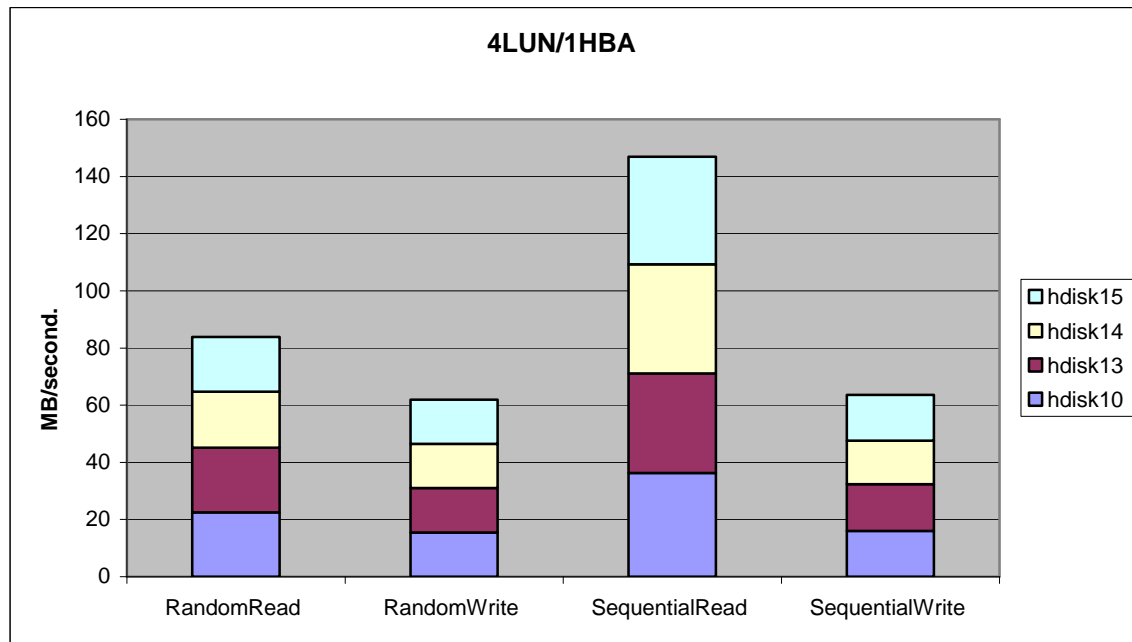
ThesecondrunistwoLUNsoveroneHBA.Ideallywewouldexpecttoseeatwofold increasein throughputforeachtest.



*Run 2*

WeseethatnotonlyrandomreadsdropperLUN,buttheaggregateisevenlower.This is alsotrueforrandomwrites,sequentialreads,andsequentialwrites.

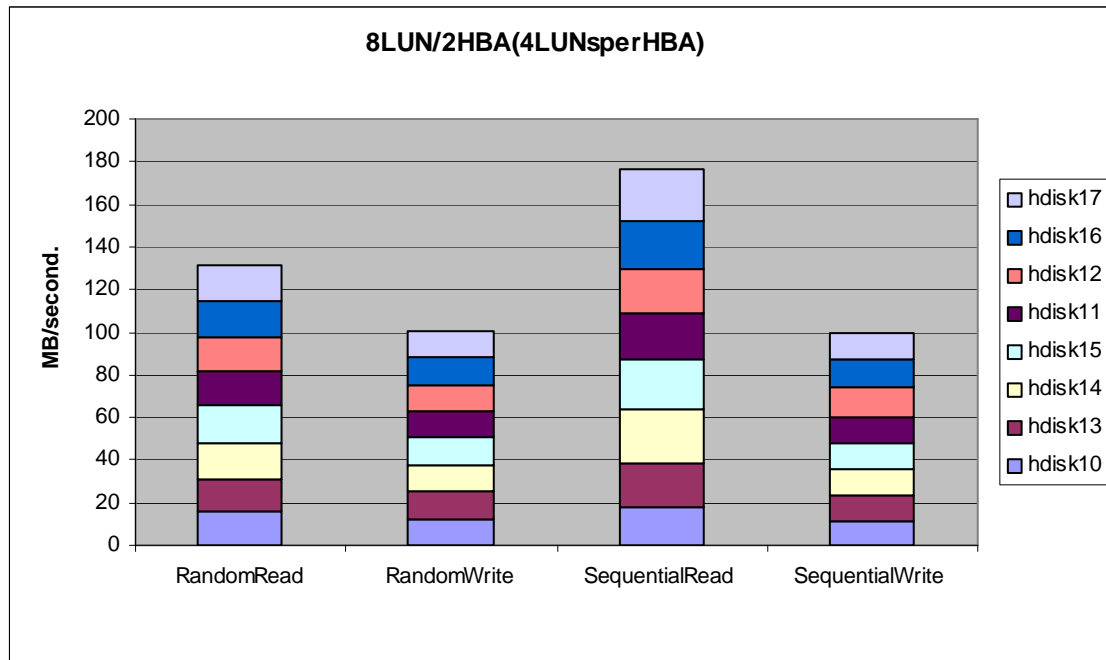
ThethirdrunisfourLUNsove roneHBA.Ideally,wewouldobserveperformance numbersbeeitheramultipleofthenumberofLUNsoverthebaselinetestsorbenearthe theoreticallimitof2Gb/secfiberspeeds(whichissomewherearound200MB/sec), whicheverislower.

*Run 3*

Random reads were observed to be an aggregated 83.8 MB/sec. This is about a 62% drop from what we might expect to see if it were a linear improvement with 4 LUNs (4 \* 54.4 or 2 Gb/sec speeds). Random writes came in at nearly 62 MB/sec. This also shows about the same 60% drop in what might be expected linearly. At this point we can see the caveat about serial ATA disk in the numbers: they are not positioned to be high I/O per second capable disk subsystems. With roughly four times the number of I/O requests coming into the controller, we see 60% loss of single LUN performance (but still achieving a higher aggregated number — unlike Run 2).

Sequential reads fell off a bit to 147 MB/sec from the single test of 181 MB/sec (a 20% drop in overall performance). Again I suspect the juggling of four request streams has something to do with this. Sequential writes appear to be unaffected at 63.6 MB/sec — nearly the same as a single sequential write. It appears the controller doesn't have a problem with I/O's when it's only being fed at 64 MB/sec.

The fourth run is 8 LUNs over 2 HBAs. This test is useful because we can determine where a bottleneck is between the controller and host port cards on the controller (fiber runs).

*Run 4*

Randomreadsperformedatanaggregaterateof131MB/sec. Thisistheculminationof twoHBAsperformingat65.6and65.4MB/secrespectively. Weseegainthereduction ofexpectedperformanceoftheprevious test –a21%dropoverthe166MB/secwe’d expecttoseebydoublingourmultipleLUN/singleHBAtest. However, theaggregate 131MB/secisthemostwe’veseenfortherandomreadtest.

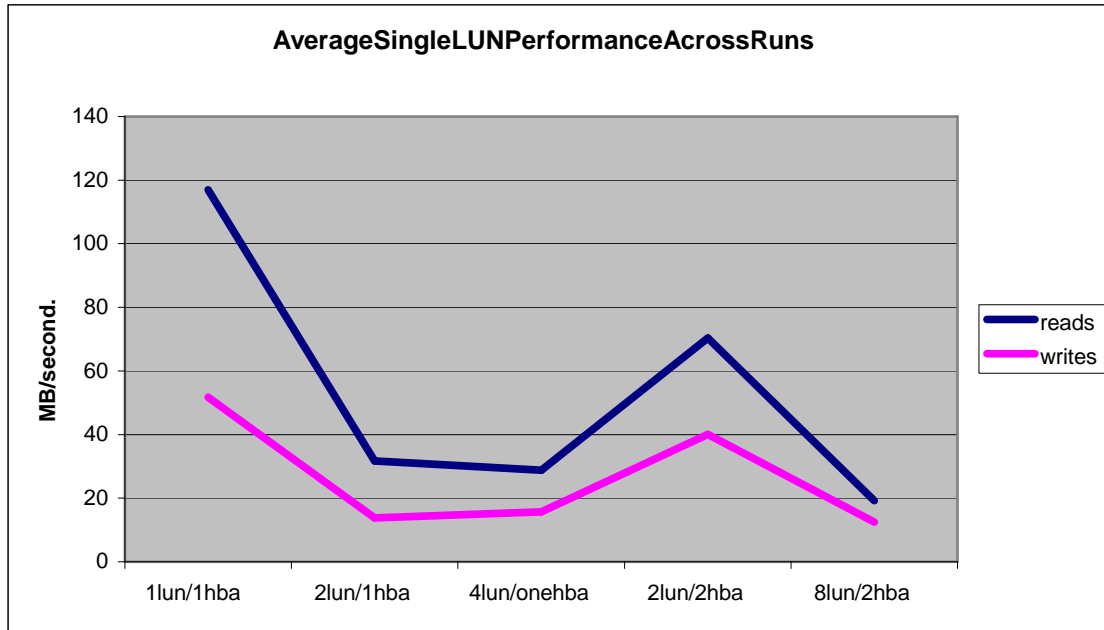
Randomwritesperformedat100MB/sec. ThisistheculminationoftwoHBAs performingat50.3MB /seceach. Thisalsoshowsasimilar18%dropinexpected performancewe’dliketoseeoverthemultipleLUN/singleHBAtestof124MB/sec(62 \*2).

SequentialReadscamein176.75MB/sec. Thisisbetterthanthesequentialreadsofthe multipleLUN/s ingleHBAtest, butit’snearlythesame(aggregated)asthesingleLUN/ singleHBAtest. Thiswouldseemtosuggestthatovertwohostports, thelimitis somewherearound180MB/seconreads. Thisisinlinewithmanufacturerclaimsof 400MB/sec, ifyou assumethatthetwounusedportswouldalsobeabletorunatthat speed. However, ifonehostportcanperformat146MB/seconmultipleLUNsequential reads, itseemslogicalthattwohostportsshouldhavebeenabletorunat292MB/sec (2x)sinceit’s stillunderthe400MB/secclaimedthroughput. Sowemustnotassumethe abilityof400MB/secforanykindoftestfromthiscontroller/diskconfigurationcombo. Wewerelimitedbytheuseof2drawersofdisk, ratherthanthe4weinitiallyhopedfor. With4drawers, wewouldhavetheuseofmorebackendchannels.

SequentialWriteswere99.75MB/sec. Againweseethe21%dropinexpected performance(bydoublingtheprevious test’s63.6MB/sec).

## Summation

Let's take a look at average LUN performance for each of the tests performed.

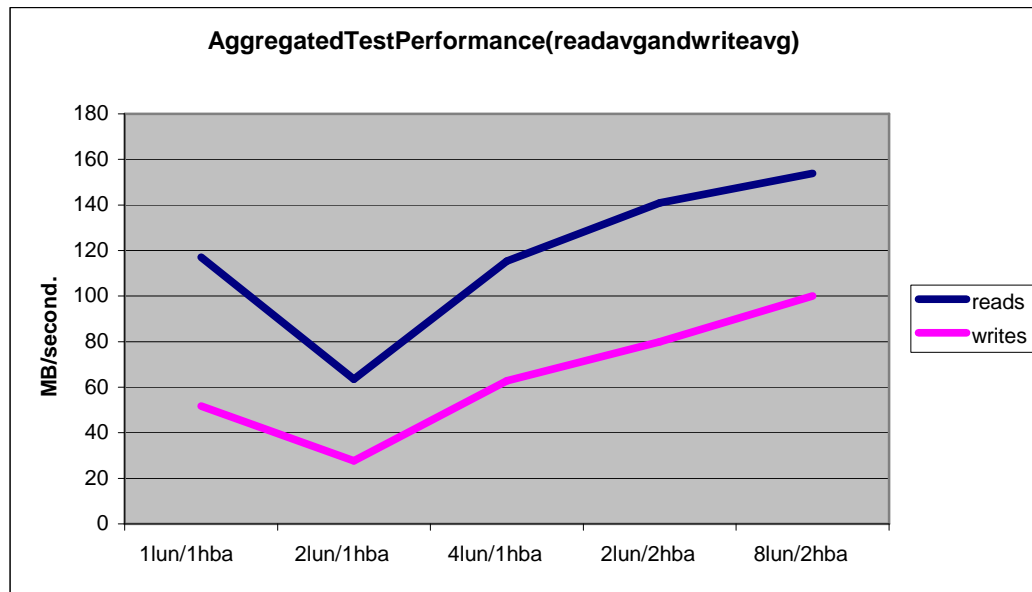


An additional test was run for this chart – 2 LUNs over 2 HBAs. You would expect to see a 2x improvement over 1 LUN/1 HBA, again this proved to be not the case.

Reads and writes were averaged together for simplicity. That is, these sequential read values and random read values were averaged together to come up with one number for each test. It's apparent, as one studies this graph, that any kind of multiple request on either read or write causes performance to drop. For use in HPSS, we might not want more than one request stream per controller.

And finally, here is a chart for the aggregated performance for each run.





Again reads were averaged and writes were averaged. This chart shows an upward trend as you add fiber channel cards. I'd like to see a four drawer Blade store configuration utilizing all four backend and four front -end ports. Beyond what is represented here is speculation.

Other less tangible things such as resistance to corruption, ease of maintenance, availability, and others were not tested due to time constraints. Testing was limited to performance runs in order to determine real -world MB/second numbers.

The point to take away from this suite of tests is that for an application that does not need more than one stream and performs sequential I/O (with emphasis on reading over writing), this disk subsystem may represent significant bang -for-the-buck. However in multiple stream scenarios where high performance I/O numbers are expected (line speed of 2Gb/sec, for example), this system falls way short of the mark.